

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4249267号
(P4249267)

(45) 発行日 平成21年4月2日(2009.4.2)

(24) 登録日 平成21年1月23日(2009.1.23)

(51) Int. Cl. F I
G06F 12/00 (2006.01) G06F 12/00 501H
 G06F 12/00 501J

請求項の数 25 (全 22 頁)

| | |
|---|---|
| <p>(21) 出願番号 特願平10-538908 (86) (22) 出願日 平成10年3月6日(1998.3.6) (65) 公表番号 特表2001-506387(P2001-506387A) (43) 公表日 平成13年5月15日(2001.5.15) (86) 国際出願番号 PCT/US1998/004567 (87) 国際公開番号 W01998/039769 (87) 国際公開日 平成10年9月11日(1998.9.11) 審査請求日 平成17年1月13日(2005.1.13) (31) 優先権主張番号 08/813,621 (32) 優先日 平成9年3月7日(1997.3.7) (33) 優先権主張国 米国(US)</p> | <p>(73) 特許権者 マイクロソフト コーポレイション アメリカ合衆国 98052 ワシントン 州 レッドモンド ワン マイクロソフト ウェイ (番地なし) (74) 代理人 弁理士 谷 義一 (74) 代理人 弁理士 阿部 和夫 (72) 発明者 ショーロフ, スリカンス, アメリカ合衆国 98029 ワシントン 州 アイザッカー サウスイースト 42 エヌディー ストリート 25049</p> |
|---|---|

最終頁に続く

(54) 【発明の名称】 ファイル・システムにおけるディスク・スペースの解放

(57) 【特許請求の範囲】

【請求項1】

持続的記憶媒体内のファイルに割り振られた記憶スペースを解放する方法であって、ファイル・システムは、前記ファイルを前記記憶媒体内の物理記憶スペース・ロケーションと関連づけるマッピング・データを維持する方法において、

前記ファイルと、ファイルの論理的先頭または論理的末尾を含んでいない、前記ファイル内の選択された論理データ・ブロックとを特定する情報を前記ファイル・システムに与えるステップと、

前記情報を前記ファイル・システムで受け取り、前記マッピング・データにアクセスして、前記選択された論理データ・ブロックに対応する、前記持続的記憶媒体内の少なくとも1つの物理ロケーションを決定するステップと、

前記少なくとも1つの物理ロケーションがもはや前記ファイルに割り振られていないことを示すように前記ファイル・システム内の前記マッピング・データを修正するステップとを備えることを特徴とする方法。

【請求項2】

請求項1に記載の方法において、前記少なくとも1つの物理ロケーションは、そのアイデンティティに対応する第1インジケータを前記マッピング・データ内にもち、前記ファイル・システム内の前記マッピング・データを修正する前記ステップは、前記マッピング・データ内において前記第1インジケータに代えて第2インジケータを書き込み、前記ファイルと関連づけて前記第2インジケータを維持するステップを備えることを特徴とする方

10

20

法。

【請求項 3】

請求項 1 に記載の方法において、前記ファイル・システムは、どのファイルにも割り振られていない空きスペースである物理ロケーションを表す空きスペース情報を維持し、前記方法は、前記少なくとも 1 つの物理ロケーションが空きスペースであることを示すように前記空きスペース情報を修正するステップをさらに備えることを特徴とする方法。

【請求項 4】

請求項 3 に記載の方法において、前記空きスペース情報は、ディスク上の前記物理ロケーションに対応するビットをもつビットマップに維持され、前記空きスペース情報を修正する前記ステップは、前記少なくとも 1 つの物理ロケーションに対応する少なくとも 1 つのビットの値を変更するステップを備えることを特徴とする方法。

10

【請求項 5】

請求項 1 に記載の方法において、前記選択された論理データ・ブロックの固定サイズを決定するステップをさらに備えることを特徴とする方法。

【請求項 6】

請求項 1 に記載の方法において、情報を前記ファイル・システムに与える前記ステップは、アプリケーション・プログラミング・インタフェースをコールするステップを含むことを特徴とする方法。

【請求項 7】

請求項 1 に記載の方法において、前記ファイルを特定する情報を前記ファイル・システムに与える前記ステップは、ファイル・ハンドルを前記ファイル・システムに与えるステップを含むことを特徴とする方法。

20

【請求項 8】

請求項 1 に記載の方法において、論理データ・ブロックを特定する情報を前記ファイル・システムに与える前記ステップは、前記ブロックのサイズを表すサイズ値を与えるステップを含むことを特徴とする方法。

【請求項 9】

請求項 8 に記載の方法において、前記サイズ値はバイト数であり、前記情報を前記ファイル・システムで受け取り、前記マッピング・データにアクセスする前記ステップは、前記バイト数をアロケーション単位数に変換するステップを含むことを特徴とする方法。

30

【請求項 10】

請求項 1 に記載の方法において、論理データ・ブロックを特定する情報を前記ファイル・システムに与える前記ステップは、前記ファイル内の前記ブロックの論理ロケーションを表すオフセット値を与えるステップを含むことを特徴とする方法。

【請求項 11】

請求項 10 に記載の方法において、前記オフセット値はバイト数であり、前記情報を前記ファイル・システムで受け取り、前記マッピング・データにアクセスする前記ステップは、前記論理ロケーションを特定のアロケーション単位に変換するステップを含むことを特徴とする方法。

【請求項 12】

請求項 1 に記載の方法において、前記ファイルの前記マッピング・データは、レコードのマスターファイル・テーブル内の前記ファイルに対応するレコード内に維持されることを特徴とする方法。

40

【請求項 13】

請求項 1 に記載の方法において、アロケーション単位はクラスタであり、前記ファイルの前記マッピング・データはクラスタ・ランのエクステンツ・リスト内に維持されることを特徴とする方法。

【請求項 14】

請求項 1 に記載の方法において、アロケーション単位はクラスタであり、前記マッピング・データはクラスタのブロックリスト内に維持されることを特徴とする方法。

50

【請求項 15】

データをストアするための複数の物理ロケーションをもつ永続的記憶媒体を含むコンピュータ・システムであって、前記コンピュータ・システムは、論理データのファイルに割り振られた物理記憶スペースを解放するように構成され、

マッピング・データをもつファイル・システムであって、該マッピング・データは前記ファイルの前記論理データを前記記憶媒体の物理ロケーションと関連づけるファイル・システムと、

ファイルの論理的先頭または論理的末尾を含んでいない、前記ファイルに割り振られた論理データのセクションを解放する要求を受け取る手段と、

前記要求を処理する手段であって、前記マッピング・データにアクセスして前記論理データのセクションに対応する少なくとも1つの物理ロケーションを決定する手段と、該少なくとも1つの物理ロケーションを前記ファイルの前記論理データから切り離すように前記マッピング・データを修正する手段とを含む手段と

を備えることを特徴とするコンピュータ・システム。

10

【請求項 16】

請求項 15 に記載のコンピュータ・システムにおいて、前記少なくとも1つの物理ロケーションはそのアイデンティティに対応する第1インジケータを前記マッピング・データの中にもち、前記マッピング・データを修正する前記手段は第1インジケータを第2インジケータでオーバーライトする手段を含むことを特徴とするコンピュータ・システム。

【請求項 17】

請求項 15 に記載のコンピュータ・システムにおいて、前記ファイル・システムはどのファイルにも割り振られていない空きスペースである物理ロケーションを表す空きスペース情報を維持し、前記要求を処理する前記手段は前記少なくとも1つの物理ロケーションが空きスペースであることを示すように前記空きスペース情報を修正する手段を含むことを特徴とするコンピュータ・システム。

20

【請求項 18】

請求項 17 に記載のコンピュータ・システムにおいて、前記記憶媒体はディスクであり、前記空きスペース情報は、前記ディスクに関連してストアされるビットマップに維持されることを特徴とするコンピュータ・システム。

【請求項 19】

請求項 15 に記載のコンピュータ・システムにおいて、前記記憶媒体の前記物理ロケーションはディスク上にクラスタを備え、前記論理データのセクションは少なくとも1つのクラスタに対応することを特徴とするコンピュータ・システム。

30

【請求項 20】

請求項 15 に記載のコンピュータ・システムにおいて、前記論理データのセクションは、バイトで表したオフセット値とそのあとに続くバイト数とで特定されることを特徴とするコンピュータ・システム。

【請求項 21】

請求項 15 に記載のコンピュータ・システムにおいて、前記ファイルの前記マッピング・データはレコードのマス・ファイル・テーブル内の前記ファイルに対応するレコード内に維持されることを特徴とするコンピュータ・システム。

40

【請求項 22】

請求項 15 に記載のコンピュータ・システムにおいて、アロケーション単位はクラスタであり、前記ファイルの前記マッピング・データはクラスタ・ランのエクステンツ・リスト内に維持されることを特徴とするコンピュータ・システム。

【請求項 23】

請求項 15 に記載のコンピュータ・システムにおいて、アロケーション単位はクラスタであり、前記マッピング・データはクラスタのブロック・リスト内に維持されることを特徴とするコンピュータ・システム。

【請求項 24】

50

請求項 1 5 に記載の コンピュータ・システム において、前記要求を処理する前記手段はアプリケーション・プログラミング・インタフェースを含むことを特徴とする コンピュータ・システム。

【請求項 2 5】

請求項 2 4 に記載の コンピュータ・システム において、前記アプリケーション・プログラミング・インタフェースは、複数のソース・ファイルを少なくとも 1 つのターゲット・ファイルに順次にマージするアプリケーション・プログラムによってコールされることを特徴とする コンピュータ・システム。

【発明の詳細な説明】

発明の分野

本発明は一般的にはファイル・システムに関し、さらに具体的には、ファイル・システムによって割り振られたディスク・スペースを解放する改良方法およびメカニズムに関する。

発明の背景

ファイル・システムが受け持つ仕事の 1 つは、ファイル内の論理データと、そこにデータがストアされている永続的記憶ボリューム (permanent storage volume) 上に置かれている物理アロケーション (割振り) 単位 (例えば、クラスタ) との間の関係をマッピング (対応づけること) することである。ファイル内の有用データ量は、そのサイズが縮小されると、そのファイルを取り扱うアプリケーション・プログラムはファイル・サイズが縮小したことをファイル・システムに通知し、ファイルに割り振られたディスク・スペースの一部を解放して再使用できるようにする。解放するデータがファイルの前方にあるときは、残余データをファイルの先頭に移し、ファイルの前方を基準にした新しいファイル・サイズをファイル・システムに知らせるのは、アプリケーション・プログラムの責任である。ファイル・システムはファイルの末尾にマッピングされたクラスタを空きスペースに戻すことで、基本的には不必要な内容をファイルの末尾から削除することによってスペースを解放する。

しかし、多くのアプリケーションは順次に、つまり、前方から後方へ向かう順序でデータを処理していく。例えば、マージ・アプリケーションでは、2 つまたはそれ以上のソート・ファイルは単一のソート・ターゲット・ファイルにマージされ、その時点でソース・ファイルは必要でなくなる。このようなマージは、ソース・ファイルの各々からデータを順次に処理していき、そのデータを該当のソート順序に従って結合し、その結合データをソート・ターゲット・ファイルに書き込むことによって行われている。マージ・アプリケーション・プログラムは大きなファイル (例えば、500 メガバイト) をマージすることがよくあるので、ソース・ファイル読取りとターゲット・ファイル書込みは、ソース・データのすべてが処理されるまでデータについて少量ずつ繰り返し行われている。ソース・データの処理が完了すると、ソース・ファイルは削除されるのが一般である。

上述したマージ方法は非常に簡易化されているが、その操作期間に大量の空きディスク・スペースが残っていることが必要である。例えば、ソース・ファイルの結合サイズの総計が 500 メガバイトであれば、ターゲット・ファイルもサイズが 500 メガバイトになる可能性がある。マージを行うためには、最大 500 メガバイトのディスク・スペースの空きがなければ、ソース・ファイルを削除することができない。このことは、ソース・ファイルがいったん削除されたあとも、総占有ディスク・スペースは未変更になっているのが一般であるので、空きディスク・スペースが基本的に一時的なものであっても同じである。もちろん、ターゲット・ファイルは一部の重複データが除去されていれば、ソース・ファイルより小さくなる場合がある。以上から理解されるように、このような大量の空きスペースはどのディスク・ボリューム上でも常に残っていると限らない。さらに、アプリケーション・プログラムに大量のデータを定期的に各ファイルの前方に移させて、そのデータが消費される時ソース・ファイルを後方から縮小できるようにすることは非常に非効率である。

上述した一時スペースの問題を解決するために、ファイルをマージする別の方法は、一緒

10

20

30

40

50

にすると大きな論理ファイルを構成する、複数の小さなファイルを管理するマージ・アプリケーション・プログラムを書くことである。このアプリケーション・プログラムは、小さなファイルが大きなファイルを構成していく様子を追跡し、データの処理時に空きディスク・スペースを作るように一部のファイルを削除することを管理する。しかし、大きな論理ファイルを構成する複数のファイルを管理することは非常に複雑である。例えば、プログラムは大きなソース・ファイルを小さなファイルに分割し、各ファイルに名前を付け、各ファイル間の論理関係を維持する必要があるため、基本的にはファイル・システム内のファイル・システムの働きをすることになる。さらに、大部分のオペレーティング・システムはアプリケーションがもつことができる、同時オープン・ファイルの数を制限しているため、同時オープン・ファイルの数が多くなるとパフォーマンスが犠牲になる。オープン・ファイルが余りに多くなるのを防止するには、アプリケーション・プログラムがさらに複雑化することになる。

10

前方から後方への順序でデータを同じように処理する他のアプリケーションとしては、先入れ先出し (FIFO) キューのファイルを取り扱うアプリケーションがある。この種のキューでは、新しいアイテムはキューの最後に追加され、必要でなくなったアイテムはキューの前方から除去される。従って、FIFO キューは、新しいアイテムをキューの最後に追加する EnQueue オペレーションと、キューが空でなければキューの前方からアイテムを除去する DeQueue オペレーションをサポートしている。IsEmpty オペレーションも用意されているが、これはキューが空であるかどうかをテストするものである。

個々のアイテムをデキュー (Dequeuing: キューから除くこと) することは持続的 FIFO キュー (persistent FIFO queue)、つまり、ディスクなどの永続的記憶媒体上にストアされた FIFO キューでは高価になる。この費用が発生するのは、デキューされたアイテム (以下、「デキュー・アイテム」という) をファイルから消去するために非常に多数の高価な入出力ディスク・オペレーションが必要になるためである。事実、持続的 FIFO キューでは、各アイテムがデキューされた直後に各アイテムをファイルから消去するのではなく、キューをクリーンアップするプログラムは複数のデキュー・アイテムをまず累積し (その際、これらのアイテムを記憶している)、そのあとでこれらのアイテムをバッチでファイルから消去している。デキュー・アイテムをこのようにバッチで消去にすると、複数のデキュー・オペレーションにわたるデキュー・コストが償却されることになる。

20

30

デキューされた (しかし、消去されていない) アイテムと、デキューされなかった残余アイテム (以下「非デキュー・アイテム」という) の両方を収めている持続的 FIFO キュー・ファイルをクリーンアップする方法はいくつかが知られている。最初の方法では、デキュー・データに残余データをオーバーライト (重ね書き) している。つまり、残余データをファイルの前方に移し、そのあと残余データのサイズに基づいてファイル・サイズを縮減している。これは、残余データのサイズに等しい一時ファイルを作成し、残余データをその一時ファイルにコピーし、そのあと残余データのコピーをファイルの前方からとってオリジナル・ファイルに戻すことによって行われている。一時ファイルはそのあとで削除されている。

第2の方法は第1の方法と似ているが、一時ファイルのコピーをオリジナル・ファイルに戻すのではなく、一時ファイルが新しい持続的 FIFO キュー・ファイルとなり、古い FIFO キュー・ファイルは削除されている。ファイル・システムは必要に応じて、ファイル・ヘッダ情報を新しい FIFO キュー・ファイルの名前で改名し、あるいは更新する。しかし、第1と第2の方法のどちらも、一時ディスク・スペースはデキューされなかったデータのサイズと等しいものを作る必要がある。さらに、第1と第2の方法では、潜在的に大量のデータがコピーされることになり、データをコピーすることは非常に費用がかかっている。

40

第3の方法では、非デキュー・データをファイル自体の中でファイルの前方に移すことによってデキュー・データに非デキュー・データをオーバーライトしている。しかし、この方法によると一時空きスペースが不要になるが、データを移動するために大量のデータをコ

50

ピーすることが依然として行われている。さらに、コピーをとっている途中でシステムに障害が起こると、ファイルが不整合の状態になるおそれがある。

最後に、アイテムは最初のファイルから最後のファイルまでにわたってシリアル番号を付けた、複数の小さなファイル内に保存しておくことができる。新しいアイテムは最後のファイルに付加されていき、そのファイルが一杯になると、新しいファイルが作成され、そのファイルが最後のファイルとなるため、ファイル総数が増加していく。最初のファイル内のアイテムがすべてデキューされると、最初のファイルは削除され、そのファイルのスペースがファイル・システムに戻される。以上から理解されるように、この方法では、複雑に階層化された、余分のファイル管理ソフトウェアを開発し、保守する必要がある。

発明の目的と概要

従って、本発明の目的は、一般的には、ファイル・システムにおいてあるファイルの任意の論理部分に割り振られたディスク・スペースを解放する方法およびメカニズムを提供することである。

さらに具体的には、本発明の目的は、ディスク・スペースの解放を高速化すると共に、大量のデータをコピーしたり、大量の一時スペースを割り振ったりしないで済むようにする、かかる方法およびメカニズムを提供することである。

本発明の別の目的は、ファイル・システムに組み込まれている、上記種類に属する前記方法およびメカニズムを提供することである。

本発明の関連目的は、システム障害から保護するためにファイル・システムが具備する既存のセーフガードと共に機能する方法およびメカニズムを提供することである。

上記目的を達成するために、本発明の関連目的はマージ・アプリケーションや持続的 F I F O キューを取り扱うアプリケーションなどの、アプリケーション・プログラムで容易に利用することができる、上記特徴を具備した方法およびメカニズムを提供することである。

本発明のさらに別の目的は、拡張可能で、ほとんどどのファイル・システムでも有効に働く、高速で、単純な、信頼性のある方法およびメカニズムを提供することである。

以上を要約すると、本発明は持続的記憶媒体内のファイルに割り振られた記憶スペースを解放する方法およびメカニズムを提供している。ファイル・システムはファイルを記憶媒体内の物理記憶スペース・ロケーションと関連づけているマッピング・データを維持している。アプリケーション・プログラムなどは、ファイルおよびそのファイル内で選択された論理データ・ブロックを特定する情報をファイル・システムに与える。ファイル・システムはこの情報を受け取ると、マッピング・データにアクセスし、選択された論理データ・ブロックに対応する、永続的記憶媒体内の物理ロケーションを決定する。ファイル・システムは、その物理ロケーションがもはやそのファイルに割り振られていないことを示すようにマッピング・データを修正し、その物理ロケーションを空きスペースに追加する。本発明の方法およびメカニズムをマージ・アプリケーションで使用すると、ソース・ファイルのデータが大きなターゲット・ファイルに追加されたときソース・ファイルを縮小することも、あるいは持続的 F I F O キューからスペースをデキューすることもできる。

本発明のその他の目的と利点は添付図面を参照して、以下の詳細な説明の中で明らかにする。

【図面の簡単な説明】

図 1 は、本発明を組み込むことができるコンピュータ・システムを示すブロック図である。

図 2 は、ファイル・システムのメタデータをディスク・ボリューム上にストアするためのテーブルを示す図である。

図 3 と図 4 はそれぞれ、本発明の一側面に従ってディスク・スペースが解放される前と後のエクステント・リストを示す図である。

図 5 は、本発明の一側面に従って修正される途中にある図 3 のエクステント・リストを示す図である。

図 6 は、修正された後の図 5 のエクステント・リストを示す図である。

10

20

30

40

50

図 7 は、本発明の一側面に従って修正された後の図 2 のテーブルを示す図である。

図 8 は、修正された後の図 3 のエクステント・リストを示す図である。

図 9 は、最適化された後の図 8 のエクステント・リストを示す図である。

図 10 と図 11 はそれぞれ、代替ファイル・システムにおいてディスク・スペースが解放される前と後のブロック・リストを示す図である。

図 12 は、本発明の一側面に従ってディスク・スペースを解放するときにとられる一般的ステップを示すフロー図である。

図 13 乃至図 16 は、ソース・ファイルがターゲット・ファイルにマージされるときのシーケンスを示す図である。

図 17 は、図 13 乃至図 16 に示すようにファイルをマージするときにとられる一般的ステップを示すフロー図である。

図 18 乃至図 20 および図 21 乃至図 23 は、デキューされ、解放されるデータと共に処理される FIFO キューを示す図である。

図 24 は、図 18 乃至図 20 および図 21 乃至図 23 に示すようにデキュー・スペースを解放するときにとられる一般的ステップを示すフロー図である。

好適実施例の詳しい説明

以下では、添付図面を参照して説明する。まず、図 1 を参照して説明すると、図 1 は本発明を組み込むことができるコンピュータ・システム全体を符号 20 で示している。図示のコンピュータ・システム 20 はサーバでも、ワークステーションでも、またはその組み合わせでもよく、公知のようにコンピュータをベースとする、1 つまたは 2 つ以上の他のリ
 ソース（資源）に接続することができる。もちろん、以下の説明で明らかにされるように、本発明はどの特定タイプのコンピュータまたはネットワーク・アーキテクチャにも限定されず、スタンドアロン型パーソナル・コンピュータなどに組み込むことが可能である。図 1 に示すように、コンピュータ・システム 20 は、オペレーティング・システム 26 がそこにロードされているメモリ 24 に接続されたプロセッサ 22 を搭載している。オペレーティング・システム 26 は Microsoft Corporation の Windows NTTM オペレーティング・システムであることが好ましい。コンピュータ 20 は、Windows NTTM ファイル・システム（NTFS）などのファイル・システム 28 をもち、これはオペレーティング・システム 26 と関連づけられているか、その内部に組み込まれている。メモリ 24 内のアプリケーション・プログラム 30 はアプリケーション・プログラミング・インタフェース（Application Programming Interface-API）を通してオペレーティング・システム 26 およびファイル・システム 28 とのインタフェースとなっている。

また、コンピュータ・システム 20 は、コンピュータ・システムを 1 つまたは 2 つ以上のネットワーク内デバイスに、キーボードおよび / またはマウスなどの 1 つまたは 2 つ以上の入力デバイス 36 に、および / またはモニタおよび / またはスピーカなどの 1 つまたは 2 つ以上の出力デバイス 38 に接続するための入出力（I/O）回路 34 も実装している。さらに、コンピュータ・システム 20 はハードディスク・ドライブ 40 などの永続的記憶媒体も実装している。

ファイル・システム 28 はハードディスク 40 上のファイルを管理し、（1）ファイルをストア、参照、共有および安全保護し、（2）ファイル・データにアクセスし、（3）ファイル健全性を維持するための、メソッドを収容しているのが一般的である。しかしながら、ファイル・システムとその関連のオペレーティング・システムは、特にこれらのファイル・システムがオペレーティング・システムに組み込まれている場合には、常に明確に区別されているとは限らない。従って、以下でファイル・システム 28 に属するとしているプロセスまたはステップはいずれも、あるいはすべてはオペレーティング・システム 26 でも実行することが可能であるものとし、その逆も同じである。

周知のように、ディスク・ドライブ 40 は複数のファイル $42_1 - 42_n$ をストアしており、これらのファイルには、論理編成のデータがディスク・ドライブ 40 上の各所ロケーションに物理的に分散されたアロケーション単位にストアされている。アロケーション単位とは任意のディスク・ボリューム上の基本的記憶単位であり、1 つまたは 2 つ以上のディ

10

20

30

40

50

スク・セクタからなるクラスタである場合がある。分散ファイル部分を論理的に連続するデータ・ブロックとして維持するために、MS-DOSのFAT (File Allocation Table-ファイル・アロケーション・テーブル) システム、Microsoftのオブジェクト・ファイル・システム (Object File System-OFS) およびNTFSなどのファイル・システム28は各ファイルのメタデータ (各ファイルに関連するアロケーション単位を順序づけている) と関連づけられたマップをストアしている。

図2は、ファイル $42_1 - 42_n$ がNTFSファイル・システムでどのような編成になっているかを示す概要図であるが、これについては、文献「Windows NTファイル・システムの内部 (Inside the Windows NT File System)」H. Custer著、Microsoft Press, 1994に記載されている。NTFSでは、アロケーション単位はクラスタであり、任意のNTFS 10
ボリュームに関する情報はレコードの集合として編成されたマスタ・ファイル・テーブル50に収められている。例えば、マスタ・ファイル・テーブルはボリュームの各クラスタごとに1ビットを持つビットマップ・レコード52を収めており、各ビットはクラスタがファイルに割り振られているか、空きスペースであるかを示す値を持っている。さらに、NTFSボリューム上の各ファイル $42_1 - 42_n$ に関する情報はレコード $54_1 - 54_n$ に収められている。他方、各レコード $54_1 - 54_n$ はデータ・フィールド $56_1 - 56_n$ を含む複数のフィールドから構成され、そこにはクラスタ・マッピング情報が入っている。しかし、NTFSでは、ファイルを構成する個々のクラスタの順序づけリストをメタデータにストアさせるのではなく、各ファイル $42_1 - 42_n$ のクラスタの連続するラン (run) 20
を追跡している。このようにすると、メタデータのスペースが節減されるが、これは、大きなファイル内のクラスタの数がそれ自体大きな数であり、NTFSはファイルのクラスタを可能な限り連続させることを試みるのが一般であるからである。

図2に示すように、 42_1 のような、ファイル内の連続クラスタの各ランごとに、データ・フィールド 56_1 はファイルの仮想クラスタ番号を表している第1の数、ディスク上のクラスタの物理ロケーションを表している第2の数、およびそのランに連続クラスタがいくつあるかを表している第3の数を含んでいる。例えば、図2に示すように、レコード 54_1 内のデータ・フィールド 56_1 は、実クラスタ10にマッピング (対応づけること) された仮想クラスタ0から始まって、11個の仮想クラスタをファイル 56_1 にマッピングしている。データは5個の連続するデータ・クラスタにわたっており、仮想クラスタ0 - 4は実クラスタ10 - 14に対応づけられている。仮想クラスタ5はファイル・データ 30
を収めている次のクラスタであり、第2のランに示すように実クラスタ19にマッピングされている。この第2のランは2の長さにはわたっているので、仮想クラスタ6はこのランによって実効的に実クラスタ20にマッピングされている。データ・フィールド 56_1 内の他のエントリをたどっていくと、レコード 54_1 に保存されている11個の物理クラスタがファイル 42_1 内の論理編成データにどのようにマッピングされているかが理解される。

クラスタ・ランはファイルのエクステント・リストと呼ばれるものを実効的に構成しており、このリストは図3に示すように、少なくとも各ラン内の最初の実クラスタ番号と、任意のファイル内の各不連続クラスタごとのクラスタ・ランの長さとを収めている。以上から理解されるように、かかるエクステント・リスト58のすべてまたは一部は各オープン・ファイル (ファイル 42_1 など) ごとにメモリ24内のファイル・バッファに置いておくと、ファイル・システム28はクラスタ・マッピング情報に高速アクセスすることが可能になる。これとは別に、ファイル・システム28がマスタ・ファイル・テーブル50のデータ・フィールド52内のメタデータを直接に利用して、必要とするファイル情報を得ようとする、エクステント・リスト58は余分なものになる。

しかし、以下では、説明を簡単にするために、エクステント・リストを中心にして本発明を説明することにする。例えば、ファイル 42_1 に、従ってデータ・フィールド 56_1 に対応するエクステント・リスト58 (図3) は同じように、ファイル 42_1 が4つのランの実クラスタ・ロケーションから構成され、これらは5の長さにはわたって (つまり、10から14までにわたって) クラスタ10から始まり、2の長さにはわたってクラスタ19に続 40
50

き、1の長さにならってクラスタ200に続き、クラスタ50から3つのクラスタ(クラスタ50から52まで)のランで終わっていることを示している。

本発明の一側面によれば、以下で詳しく説明するように、アプリケーション・プログラム30はファイルのデータの任意のセクションを解放、つまり、デコミット(decommit)することができ、ファイルの末尾だけでスペースを解放することに限定されない。スペースの解放を行うには、アプリケーション・プログラム30はファイル・システム28にコールを出し(API 32内の定義されたAPI 46から)、必要でなくなったが、まだファイルに割り振られたままのデータのセクションをデコミットする。好ましくは、データのセクションは1つまたは2つ以上のクラスタに対応する事前定義サイズの倍数になっている。

10

スペースを解放するには、ファイル・システム28はアプリケーション・プログラムからほとんど見えないような形でファイルのメタデータを操作する。具体的には、アプリケーション・プログラム30は、デコミットされたスペースがファイルの一部としてまだ残っているものとして、その内部ポインタとオフセットを維持している。もちろん、デコミットされたデータ・ブロックがファイルに属さなくなると、そのファイルに読み書きできなくなるが、アプリケーション・プログラム30はスペースのデコミット前にそのデータまたはスペースが必要でなくなったことを確かめること以外は、なにも調整を行う必要はない。以下で詳しく説明するように、アプリケーション・プログラムがデコミットされたスペースの読み取りを試みると、ゼロが戻されることがあり、書き込みを試みると、デコミットされたスペースが再度デコミットされることがある。

20

次に、本発明の動作について説明すると、本発明のAPI 46がまずアプリケーション・プログラム30からコールされるとき、プログラム30はデコミットすべきスペースを含んでいるファイルのハンドル(例えば、ハンドルが5であるファイル42₁)を含む情報をAPI 46に渡す。アプリケーション・プログラム30からは、デコミットを開始する論理ロケーションまでのオフセット(例えば、ファイルの先頭からのバイト数を示すオフセット・ポインタ)とデコミットする長さの値(つまり、バイト数などのスペース量)も渡される。例えば、ファイル42₁のコールでは、5に等しいファイル・ハンドル、1メガバイトのオフセットおよび512キロバイトの長さを渡すことができる。好適実施例では、ファイル・ハンドルはオープン・ファイルに対応し、それ以外はAPI 46からエラー・メッセージが戻される。

30

単純性を保つために、バイト数で表したオフセットと長さの値は、好ましくはファイル・セクションと呼ばれる単位の倍数になっている。一般的に、ファイル・セクションのサイズは任意であるが、ファイル・システム28内では固定されており、ファイル・システム28がどのようにデータを編成するかと若干の関係をもって選択されている。例えば、NTFSでは、ファイル・セクションのサイズはディスク・ボリューム上のクラスタのサイズの倍数であり、試みとして64キロバイトが選択されている。この64キロバイトのサイズは512バイト・クラスタ、1キロバイト・クラスタ、2キロバイト・クラスタ、以下同様の倍数であり、これらはいずれもNTFSで使用可能である。OFSでは、ファイル・セクション・サイズは256キロバイトである。このサイズも、ディスク40の入出力オペレーション数が、各オペレーションで解放されるスペース量とバランスがとられるように選択されている。つまり、多数の小さなデコミットを行わないで済むだけ十分に大きい、デコミットを行う前に大量の一時ディスク・スペースが必要でないくらい十分に小さくなっている。もちろん、以下で明らかにするように、アプリケーション・プログラム30は1回のAPIコールで2つ以上のファイル・セクションをデコミットすることが可能である。例えば、1回のAPIコールで、8個の64キロバイト・ファイル・セクションに相当する512キロバイトをAPI 46にデコミットさせることが可能である。アプリケーション・プログラム30がAPI 46をコールしたとき、オフセットおよび/または長さがファイル・セクションの倍数になっていないと、API 46からエラー・メッセージなどが戻される。しかし、適切なアプリケーション・プログラム30は、ファイル・セクションのサイズが事前に知らされているか、あるいはファイル・システム2

40

50

8に問い合わせたファイル・セクションのサイズを知ることができるようになっている。従って、別の方法として、オフセットおよび長さの値の単位としてバイト数ではなく、ファイル・セクションをアプリケーション・プログラム30が渡すようにすることも可能である。

上記とは別に、より複雑ではあるが、ファイル・セクションのサイズに対応する値ではなく、そのオフセットと長さパラメータの任意のバイト値をアプリケーション・プログラム30が渡すようにすることも可能である。このような代替方法では、ファイル・システム28はオフセットと長さをファイル・セクションに変換し、可能な限り多数のファイル・セクション（ゼロであることも可能）をデコミットする。次に、ファイル・システム28はデコミットされなかったバイト数に相当する残余をアプリケーション・プログラム30に戻すか、さもなければデコミットのエクステントを示すポインタを戻す。その後、アプリケーション・プログラム30はデコミットするポインタを戻り値に基づいて調整することができる。別の方法として、ファイル・システム28が自身で残余を追跡し、アプリケーション・プログラム30が部分的ファイル・セクションの残余をデコミットすることをあとで要求したとき、そのバイトをデコミットすることも可能である。

上記とは別に、好ましい方法は、アプリケーション・プログラム30が事前定義のファイル・セクション・サイズを使用することである。従って、以下の例は好ましいAPIコール条件に基づいている。さらに、64キロバイトのような代表的なファイル・セクション・サイズを取り扱うとき、オフセット・パラメータに対応するクラスタはエクステント・リスト内に深く埋め込まれる可能性があるために、その後でデコミットするクラスタの数が相対的に大きな数になることがある。このような大きな数を扱うことを避けるために、以下の例ではその目的上、特に断りがない限り、任意のファイル・セクションのサイズがクラスタのサイズに等しいものとしている。つまり、どの特定クラスタもデコミット可能になっている。若干実用目的に合わないが、以上から理解されるように、このような小さなファイル・セクション・サイズの使用が可能であるため、以下の例は単純化されているが、そのことは本発明の精神と範囲から逸脱するものではない。

図3はファイル・ハンドルが5であるファイル42₁の場合のエクステント・リスト58を示している。第1の例では、アプリケーション・プログラム30はAPI 46をコールし、20キロバイトのオフセットから始まり、8キロバイトの長さにわたって、ハンドル5をもつファイルをデコミットすることを要求している。この例では、クラスタ・ファクタ（cluster factor）はクラスタ当たり4キロバイトであり、1つのファイル・セクションは4キロバイトのサイズになっている。API 46がコールされると、図12のステップ100でAPI 46はファイル・ハンドルを使用して、エクステント・リスト58などの該当ファイル情報を選択する。ステップ102で、API 46は、20キロバイト・オフセットをクラスタ当たり4キロバイトのクラスタ・ファクタで除することによって、そのオフセットをファイルの先頭から5クラスタ離れたオフセットに変換する。エクステント・リスト内の最初のエントリはクラスタ10、長さ5であるので、ファイル・システムは、クラスタ10、11、12、13および14が最初の5クラスタとして、この順序でファイル42に割り振られていることを知っている。従って、オフセット値は第2のランのクラスタ19を指している。ステップ104で、長さパラメータは8キロバイトをクラスタ当たり4キロバイトのクラスタ・ファクタで除することにより2クラスタに変換される。従って、図3に示すように、クラスタ19と20がデコミットされることになる。

次に、ステップ106で、デコミットするクラスタ、この例ではクラスタ19と20がラン境界上に置かれているかどうか決定される。言い換えれば、ステップ106では、クラスタ19がランの先頭にあり、クラスタ20がランの末尾にあるかどうか決定される。この例では、クラスタ19と20はランを開始し、終了しているため、このプロセスはステップ110にブランチする。

ステップ110で、クラスタ19と20はメタデータを修正することにより、具体的には、クラスタ19と20がもはやファイルの一部でないことを示すようにエクステント・リ

10

20

30

40

50

スト58の値を変更することによってデコミットされる。図4に示すように、マイナス1(-1)の値は実クラスタ番号19をオーバーライトするが、有効なクラスタ番号を表していない、どの事前定義値でも表示値になるように選択することも可能である。最後に、ステップ112で、デコミットされたクラスタ19と20は空きスペースに戻される。NTFSシステムでは、これはクラスタ19と20に対応するビットの値を、ビットマップ52内でトグルすることによって行われる。ステップ110と112は、システム障害が発生したときオールオアナッシング(all-or-nothing)オペレーションを保証するように実行されることが好ましい。

図3および図5乃至図6は、デコミットされる1つまたは複数のクラスタがラン境界に正確に一致していない例を示している。この例では、エクステント・リストが図3に示すものと同じであり、図3と同じようにクラスタ・ファクタがクラスタ当たり4KB、ファイル・セクション・サイズが4KBであるとして、クラスタ13をデコミットしようとしている。上述したように、アプリケーション・プログラム30は、5に等しいファイル・ハンドル(ファイル42₁)、12キロバイトのオフセット、および4キロバイトの長さを指定してAPI 46をコールすることによってかかるスペースの解放を開始する。同じく上述したように、これらのパラメータを使用して、ステップ100-104はオフセットと長さパラメータをクラスタ情報に変換し、クラスタ13だけがデコミットされるものと決定する。

しかし、ステップ106では、クラスタ13がランに一致していないで、もっと大きなランの一部であると決定される。その結果、ステップ106はステップ108にブランチし、そこでエクステント・リスト58内のメタデータが、図5に示すように、10長さ5のランを10長さ3、13長さ1および14長さ1の複数のランに分割することによって修正される。例えば、この修正はエクステント・リスト58のランを一時スペースにコピーし、そのランをエクステント・リスト58に戻すように再コピーし、新しい情報を挿入し、該当する場合には既存の情報を修正することによって行うことができる。なお、図5には事実上図3と同じクラスタがリストされているが、図5ではクラスタ13はラン境界と一致している。

ステップ108に続いて、スペースは、クラスタ13がファイル42₁にもはや割り振られていないことをエクステント・リスト58に示すことによってステップ110でデコミットされる。クラスタ13をデコミットしたことは図6の修正エクステント・リスト13に示されているが、そこでは-1がクラスタ13の個所に書かれている。最後に、ステップ112で、デコミットされたクラスタ13はビットマップ52を上述したように修正することによって空きスペース(NTFSの)に戻される。

以上から理解されるように、メタデータの操作はエクステント・リストを通してではなく、マスタ・ファイル・テーブル50内で直接的に行うことが可能である。図7は、図3、図5および図6を参照して上述したように、クラスタ13がデコミットされたときそれに応じて、レコード54₁のデータ・フィールド56₁が直接に修正されると、ランがどのようになるかを示している。なお、図7には図6と同じ情報が示されているが、ここには、各ランの論理クラスタ値も示されている。

図8および図9は複数のランが1回のAPIコールでデコミットされる例を示している。図8では、ゼロのオフセットから始まって、図3の最初の9クラスタは1回のAPIコールまたは一連のAPIコールで、APIコールによってすでにデコミットされている。図9では、デコミットされたクラスタは連続するデコミット・クラスタの長さを合計することによって、単一の連続デコミット・スペース・ランに結合されている。本発明では必要でないが、このような最適化を定期的に行うと、エクステント・リスト58内の(またはデータ・フィールド56₁内の)エン트리数が大きくなり過ぎるのを防止することができる。

本発明の一側面によれば、デコミットされたクラスタ・ランをエクステント・リスト58から単に除去するだけでなく、ファイル・システム26はそのランにデコミットされたとのマークを付ける。その結果、デコミットされたスペースはまだファイル44₁と論理

10

20

30

40

50

的に関係づけられているので、ファイル42₁を基準にした、アプリケーション・プログラム30のオフセットとポインタ値は有効のままになっている。例えば、アプリケーション・プログラム30がファイル42₁の先頭からの32キロバイトのオフセットをファイル・システム28に送るときはいつでも、オフセットは、デコミットされたファイル・セクションがあるかどうかに関係なく、正しい物理クラスタ（図3または図4のどちらも、クラスタ50）を指している。

しかし、アプリケーション・プログラムの内部データ・ポインタが未変更であっても、アプリケーション・プログラム30は戻って、デコミットされたスペースを読み書きすることができない。この目的のために、アプリケーション・プログラム30がファイル・システム28に読取り操作を要求したが、渡したオフセットまたは長さがデコミットされたスペースに対応するものであれば、ファイル・システム28はそのスペースがすでにデコミットされていたと、ファイルのメタデータから判断し、デコミットされたファイル・セクションのバイトに対してゼロ（さもなければエラー）を戻してくる。事実、API 46が試みとして"WriteZeros()"と名づけられているのは、アプリケーション・プログラム30は特定ファイルのデコミットされたスペースを実効的にゼロにするからである。なお、デコミットされたスペースは空きスペースになるので、他のファイルで使うことができ、そこには非ゼロのデータが入るのが代表的であるが、デコミットされたスペースはそのスペースがデコミットされたファイルでは論理的にゼロだけを収めている。一部の高度セキュア・ファイル・システムでは、空きクラスタを物理的にゼロにしている。

同様に、あるファイルのデコミットされたスペースにデータを書き込もうとすると、そのファイルではエラーとみなされる。しかし、ファイル・システム28に新しいスペース（これは偶然にデコミットされたスペースであることがある）を割り振らせ、そこにデータを書かせることも可能である。そのようにすると、ファイル・システム28は新しく割り振られたスペースのロケーションとサイズを反映するように、そのファイルのメタデータを調整すること、つまり、新しいクラスタ・ランを正しいロケーションでメタデータに挿入すること（該当する場合には既存のランを延長すること）も必要になる。

代替ファイル・システムがアロケーション単位をファイルにマッピングするもう1つの方法は、ブロック・リストによるものである。公知のように、ブロック・リストは単純にファイルに関連するアロケーション単位のリストである。例えば、FATシステムでは、ファイルの各クラスタはファイルの次のクラスタ（またはファイル末尾マーカ）を指しているエントリをファイル・アロケーション・テーブルの中にもっている。その情報から、種々データ・クラスタ間の関係をマッピングしているブロック・リストはそのファイルに関連するバッファに（または持続的メタデータ・ストレージに）置かれる。ファイル42₁の場合の例示ブロック・リスト60（これはNTFSシステムでは図3のエクステンツ・リスト58に相当する）は図10に示されている。

図10は、ディスク・ボリューム上で利用可能な空きスペースを示すブロック・リスト62も示している。このブロック・リスト62内の情報はメモリに置いておくことも、ディスク・ドライブに置いておくこともできるが、どの場合も、ファイル・システム28が利用できるようになっている。

図11は、ファイル42₁のブロック・リスト60およびファイル・ハンドル5、オフセット20KBおよび長さ8KBのパラメータを指定してWriteZeros() API 46がコールされた後の空きスペースのブロック・リスト62を示している。前述したように、図12のステップ100は該当するブロック・リストを選択し、ステップ102はオフセットをクラスタ19を指すポインタ（先頭から5クラスタ）に変換し、ステップ104はデコミットする長さを2クラスタに変換する。しかし、ステップ106（場合によっては、ステップ108も）が実行される代わりに、ブロック・リストを使用するファイル・システムでは、ステップ104は即時にステップ110にジャンプする。これは、各クラスタが個々にリストされる時（固有長さは1になっている）取り扱うべきランがブロック・リストになく、従って、デコミットが実効的に常に境界に一致しているからである。

ステップ110で、クラスタ19と20のブロック・リスト60内のエントリはそこにイ

10

20

30

40

50

ンジケータ（-1など）を入れることによってデコミットされる。ステップ112で、クラスタ19と20は図11のブロック・リスト62に示すように、空きスペースに追加される。

使用されるのがブロック・リストであるか、エクステント・リストであるかに関係なく、アプリケーション・プログラム30ではなくファイル・システム28を使用してスペースのデコミットを管理する利点は、多くのファイル・システムが、NTFSを含めて、システム障害から保護するセーフガードの働きをすることである。例えば、NTFSはオールオアナッシングのトランザクションを保証するロギング・ファイル・システムである。つまり、各オペレーションはアトミック単位で実現されている。あるトランザクションの完了前にシステム障害が起こると、NTFSはそのログを使用してオペレーションをロールバックする。

10

以下では、本発明をマージ・アプリケーション・プログラム（例えば、アプリケーション・プログラム30はマージ・アプリケーションである）で利用した場合について説明するが、その概要図は、図13乃至図16と図17のフロー図に示されている。公知のように、マージ・アプリケーション・プログラムは、2つまたはそれ以上のソート・ソース・ファイルからのデータを1つのソート・ターゲット・ファイルに結合するものである。例えば、アプリケーション・プログラムは、名前とアドレスのアルファベット順リストを同種の別のリストとマージして、ソース・リストの結合サイズと等しいアルファベット順リストを得ることができる。

マージ・アプリケーション・プログラム30の予備ステップ200（図17）では、少なくとも2つのソース・ファイル70、72（図13）がオープンされ、ターゲット・ファイル74（図13）が作成される。スペースを節減するために、以下で説明するように、ターゲット・ファイル74は最初は1ファイル・セクションだけのサイズで、例えば、64キロバイトでオープンされる。なお、前述したように、アプリケーション・プログラム30は対応するファイル・システム28のファイル・セクションのサイズを知っている。ステップ202はソース・データが読み取り可能になっているかどうかをテストして決定する。この時点ではソース・ファイル70、72はオープンしたばかりであるため（およびこの例では、長さが非ゼロであるため）、読み取るデータがあるので、ステップ202はステップ204にブランチする。ステップ204で、最初のデータ・ブロック（例えば、20キロバイト）が各ソース・ファイル70、72から読み取られ、データを公知のようにマージすることによって処理される。この例では、マージ・データはサイズが36キロバイトになる場合がある。なお、1つのソース・ファイルから処理されるデータ量は別のソース・ファイルから処理されるデータ量と同じでないのが一般である。従って、ループを通るたびに、追加データはステップ204で特定ソース・ファイルから必要になったときだけ読み取られる。

20

30

図13は、2つのソース・ファイル70、72をターゲット・ファイル74にマージする初期ステージにおけるマージ・アプリケーション・プログラム30を示し、そこでは陰影を付けたエリアはファイル内のデータを表している。単純化のために、ソース・ファイル70、72は、それぞれが4つのファイル・セクションと5つのファイル・セクションからなるものとして示されているが、ソース・ファイルはどのサイズであっても構わず、必ずしもファイル・セクションが正確に倍数であるとは限らない。アプリケーション・プログラム30は第1ソース・ファイル70の先頭（図13ではゼロ・バイト）を示す第1ポインタ76または類似手段と、ソース・ファイル70からデータが読み取られた（または下述するようにバッファに置かれ、処理され、および/または書き込まれた）個所を示す第2（読み取り）ポインタ78または類似手段とを保持している。同様に、アプリケーション・プログラム30は第2ソース・ファイル72の先頭（図13ではゼロ・バイト）を示す第3ポインタ80と、ソース・ファイル72からデータが読み取られた個所を示す第4（読み取り）ポインタ82とを保持している。アプリケーション・プログラム30はデータがターゲット・ファイルに書き込まれた個所を追跡するためのターゲット（書き込み）ポインタも保持している。

40

50

ステップ206は処理されたデータ(40キロバイト)がターゲット・ファイル74の残余サイズに収まるかどうかをチェックして決定する。この例では、ターゲット・ファイルはオープンされたばかりであるので、64キロバイト全部が使用されて処理データの36キロバイトで満たされることになる。従って、プロセスはステップ210にブランチする。十分なスペースがターゲット・ファイル74に残っていなければ、ステップ206はステップ208にブランチし、そこでターゲット・ファイル74のサイズはファイル・システムに追加ディスク・スペースを要求することによって増加されることになる。なお、必要時にだけターゲット・ファイル74にスペースを追加するようにすると、非常に大きなターゲット・ファイルを初めに作っておかなくてもマージを行うことができる。

どの場合も、ステップ210では、処理されたデータはターゲット・ファイル74に書き込まれる。以上から理解されるように、システム障害から保護するために、バイトの読み取りを追跡するのではなく、ポインタ78と82はデータが各ソース70、72から読み取られ、処理され、無事にターゲット・ファイル74に書かれた個所を追跡することができる。これにより、バッファに置かれたデータがシステム障害発生時に失われることが防止される。従って、ポインタ78、82はその時点でそれに応じて増加させ、図13乃至図16の右に移動することができる。例えば、図14に示すように、ポインタ78と82は図13のその位置に対して右に移動している。

ステップ212(図17)は、1つのファイル・セクション・サイズ(この例では、64キロバイト)より大きいデータがソース・ファイル70から読み取られたか(または処理されて無事に書かれたか)どうかを決定する。例えば、データがどれだけソース・ファイル70から読み取られたかを決定するために、プロセスはまずポインタ78とポインタ76との差を計算する。図14に示すように、ポインタ78と76との差は1ファイル・セクションより大きくなっている。従って、ステップ214で、ファイル70内で"1"で示されているファイル・セクションが解放される。

プロセスのステップ214は、上述したようにWriteZeros() API 46をコールすることによってファイル70の前方からファイル・セクションをデコミットする。この目的のために、API 46には、ファイル70のファイル・ハンドル、オフセットのバイト数(図14ではゼロ)およびデコミットする長さ(64キロバイト、つまり、1ファイル・セクション)が渡される。同じく上述したように、WriteZeros() API 46はファイル70に対応するメタデータを操作することによってこれを行う。なお、データはマージ・アプリケーションでは順次に読み取られ、処理されるので、ファイル・セクションは前方からデコミットされる。

図15は図のファイル・セクション"1"が読み取られた後のファイルの状態を示している。ポインタ76はファイル・セクション"2"の先頭まで、例えば、64キロバイトだけ進められ、ファイル・セクション"1"がデコミットされたことを示している。アプリケーション・プログラム30は、「ゼロ」がいま解放されたファイル・セクション"1"の先頭をまだ示すようにそのポインタとオフセットを保持しているが、ファイル・システム28はそのスペースに空きのマークを付けている。

ステップ214に続いて、ステップ216はソース・ファイルのすべてがテストされたかどうかを決定し、すべてがテストされていないならば、ステップ218にブランチして次に評価すべきソース・ファイルを選択する。従って、この例では、ソース・ファイル72はまだテストされていないので、ステップ216はステップ218にブランチし、そのあとステップ212が再び実行されるが、この時はファイル70のポインタ80、82が使用される。しかし、図14に示すように、およびステップ212で決定されたとおりに、ファイル72では、ポインタ82とポインタ80との差は1ファイル・セクション以下である。従って、スペースは解放されないで(今回はステップ214はスキップされる)、プロセスは直接にステップ216に戻る。テストする必要のあるソース・ファイルは残っていないので、プロセスはステップ202に戻り、追加のデータを読み取る。

図16はさらに以後のステージでのマージを示している。図16に示すように、ソース・ファイル70と72はターゲット・ファイル74の増加と共に縮小している。その結果、

10

20

30

40

50

ファイル・システム 28 は以前にソース・ファイル 70 と 72 に割り振られていたスペースをターゲット・ファイル 74 に追加することができる。以上から理解されるように、この手法によると、ファイル・システム 28 はどの時点でも、少量の一時スペースを割り振るだけで済むことになる。事実、ファイル・セクションがそのソース・ファイルから除去されたかどうかを、なんらかの方法でアプリケーション・プログラム 30 がチェックして確かめるとした場合、最悪のケースでは、総一時スペース量はソース・ファイルの数にファイル・セクションのサイズをかけたものにほぼ等しくなる。最善のケースでは、総一時サイズは 1 ファイル・セクションに若干の増加分を加えたものである。

最終的には、ファイル 70 と 72 のどちらからも読み取るべきソース・データがなくなることになる。その時がくると、ステップ 202 (図 17) はステップ 220 にブランチし、そこで残っているソース・ファイルのスペースがあれば解放されることになる。例えば、ソース・ファイル 70 と 72 が削除される。

次に、図 24 のフロー図に概要を示すように、持続的 F I F O キューで利用したときの本発明について説明する。図 18 は以前にデキューされたが、ファイル 86 から消去されていないデータ量 (セクション " A " で示されている) をもつ持続的 F I F O キュー・ファイル 86 を示している。ここではそのように示されていないが、" A " はゼロ・バイトであってもよい。ポインタ 88 または類似手段はデキューされたスペースが始まるロケーション (例えば、バイト・オフセット) を追跡し、別のデータ・ポインタ 90 または類似手段はファイル 86 のセクション " C " 内のデキューされていない残余データ・アイテムの先頭を追跡する。なお、スペースが上述したように以前にデコミットされていれば、デキューされたスペースの先頭はアプリケーション・プログラム 30 側から見たときバイト・ゼロになく、このことはデータ・ポインタ 88 内の非ゼロ値によって示される。例えば、ポインタ 88 は 450 キロバイトの位置を示すことができるのに対し、ポインタ 90 は 475 キロバイトを示すことがある。

DeQueue オペレーションが図 24 のステップ 300 で実行されると (これにより 1 つまたは 2 つ以上のアイテムがファイル 86 の前方から除去される)、別のデータ量がデキューされている。この追加のデキュー・データは図 19 にセクション " B " として示されている。以上から理解されるように、デキュー・オペレーションの一部として、ポインタ 90 は残余の非デキュー・セクション " C " データの新しい開始位置を示すように調整される。この例では、ポインタ 90 は 500 キロバイトの位置を示すように増加されている場合がある。以上から理解されるように、図 24 に示すプロセスは各アイテムがデキューされた後ではなく、いくつかのデキュー・オペレーションの後で実行することができる。

デキュー・オペレーションに続いて、ステップ 302 で、デキューされた総スペース量が、例えば、ポインタ 88 にストアされたバイト値をポインタ 90 にストアされたバイト値から差し引くことによって計算される。この総サイズは " A " と " B " の結合セクションで図 19 に示されており、この例では、50 キロバイト (つまり、500 キロバイト - 450 キロバイト) に等しくなっている。ステップ 304 で、この総デキュー量は図 18 乃至図 20 にインターバル " F S " で示されているファイル・セクションのサイズと比較される。この例では、ファイル・セクションはサイズが 64 キロバイトであり、これは固定され、アプリケーション・プログラム 30 に知らされている量である。50 キロバイトはファイル・セクションの 64 キロバイトより小であるので、この時点ではスペースはデコミットされず、ステップ 304 はブランチして図 24 のプロセスから出る。しかし、次回に図 24 のプロセスを通るときは、セクション " A " からスタートするのではなく、ポインタ 90 はセクション " A " と " B " の両方がデキューされていた個所より先に進んでいる (つまり、この例では 500 キロバイトのオフセット)。その結果、F I F O キュー・ファイル 86 は以前にデキューされたが、ファイル 86 から消去されていない新しいデータ量 (つまり、50 キロバイト) になっており、これは図 20 にセクション " D " で示されている。

図 18 乃至図 20 とは対照的に、図 21 乃至図 23 は総デキュー・スペースがファイル・セクションより大である場合を示している。以前と同じように、図 21 は以前にデキュー

10

20

30

40

50

されたが、ファイル 86 から消去されていないデータ量（ゼロ・バイトのこともある）をもつ F I F O キュー・ファイルを示している。このデータはセクション " E " で示され、例えば、700 キロバイトから始まり、730 キロバイトマイナス 1 バイトにわたっている場合がある。なお、図 21 乃至図 23 では、ポインタ 88 と 90 は図 18 乃至図 20 を参照して前述したものと類似のデータ・ポインタ情報を保持しており、この例では、それぞれ 700 キロバイトと 730 キロバイトの値をもっている（従って、その値を指している）。

図 22 に示すように、Dequeue オペレーションが図 24 のステップ 300 で実行されたあと、セクション " E " で示すように別のデータ量がデキューされている。この場合も、ポインタ 90 は残余の非デキュー・セクション " G " の新しい開始位置を示すように調整され、例えば、780 キロバイトに増加されている。しかし、図 22 に示すように、今度は総デキュー・スペース（" E " + " F "）は 1 つのファイル・セクション、つまり、この場合もインターバル " F S " で示されている図 21 乃至図 23 のファイル・セクション・サイズより大になっている。

従って、ステップ 302 で総デキュー・スペース量は、例えば、ポインタ 88（700 キロバイト）をポインタ 90（780 キロバイト）から差し引くことによって計算され、80 キロバイトが得られる。ステップ 304 でこの総デキュー・スペース量（80 キロバイト）はアプリケーション・プログラムによってファイル・セクションの既知サイズである、64 キロバイトと比較される。しかし、今回はデキュー・スペースはファイル・セクション・サイズより大であるので、ステップ 304 での比較はステップ 306 にブランチする。

本発明によれば、アプリケーション・プログラムはファイル 86 の前方からスペースをデコミットするファイル・セクションで A P I をコールする。例示のコールは WriteZeros（5, 700KB, 64KB）にすることができ、ここで " 5 " はファイル・ハンドルである。なお、持続的 F I F O キューのアプリケーション・プログラム 30 は、デキュー・スペースをファイル・セクションのサイズで除するだけで 2 つ以上のファイル・セクションをデコミットすることができ、その商はデコミットするファイル・セクションの数と等しくなる。別の方法として、プロセスは残余のデキュー・スペースが 1 つのファイル・セクション以下になるまでループバックすることによって、ファイル・セクションを 1 つずつ繰り返しデコミットすることも可能である。

ステップ 308 で、ファイル・システム 28 に関連する A P I 46 は該当数のファイル・セクションを解放し、解放したスペースを利用可能ディスク・スペースに加える。これは、前述したようにファイルのメタデータと空きスペース・ビットマップなどを操作することによって行われる。ステップ 310 で、アプリケーション・プログラム 30 はポインタ 88 を、デコミットされたなかった残余デキュー・スペースの先頭に移す。この例では、ポインタ 88 は 764 キロバイト（700 キロバイト + 64 キロバイト）に増加する。この増加を分かりやすく示したのが図 23 であり、セクション " H " はデキューされたが、ファイルから消去（デコミット）されていない残余スペースを示している。

以上から理解されるように、ファイル 86 の末尾以外から、この例では前方からスペースを解放すると、高速にかつ安価な方法で持続的 F I F O キューを消去することができる。事実、この消去操作は、データをコピーする必要もなければ、複雑なファイル操作の必要もなく、必要ならば、アイテムがデキューされるつど実行できるくらい効率的なものである。

しかし、消去操作によると物理ファイル・サイズは保たれるが、論理ファイル・サイズは持続的 F I F O キューや場合によっては他のアプリケーションによると大きくなっていく。例えば、ファイルの前方はスペースがデコミットされると縮小するとしても、論理ファイル・サイズは新しいアイテムがキューに追加されると大きくなっていく。また、ポインタはその値が大きくなっていく。これは、ファイル・サイズとポインタが論理的に 2^{64} バイトまで大きくなることのできる N T F S では問題にならないことが認められているが、他のファイル・システムでは、最終的には、この成長を続ける数を取り扱うとき問題が起

10

20

30

40

50

こる可能性がある。

その結果、本発明ではその必要はないが、このような暴走的成長の可能性を認識し、望ましいときにファイル・サイズとポインタを実効的にリセットするようにアプリケーション・プログラムを書くことができる。この目的のために、アプリケーション・プログラム30は有用(デキューされていない)データを一時ファイルにコピーし、オリジナル・ファイルが占有しているスペースが解放される。つまり、オリジナル・ファイルは削除される。その後、メタデータは一時ファイルを指すように再マッピングされ、ポインタ88、90はゼロにリセットされる。実際には、ファイルの前方にあるホールが除去される。ある種のデータ・コピーを行うと、上記のような最適化が得られるが、この最適化は頻繁に行う必要はない。例えば、このような最適化を行う望ましい時期は残余の有用データが相対的に少ないときであり、大量にコピーする必要がないからである。もう1つの望ましい時期は、ファイル・サイズまたはポインタの値がある種の、あらかじめ決めたしきい値の安全量を超えたときである。

10

その内部ポインタを調整するアプリケーション・プログラム30に関連してファイル・システム28にこの最適化を行わせることも可能である。例えば、アプリケーション・プログラム30による要求が別のAPIを通して行われたとき、デコミットされたクラスタのランをファイルのメタデータから除去し、ファイル・サイズをそれに応じて調整することができる。同様に、FATシステムなどのファイル・システムでは、最初のクラスタを指している情報および/またはFATテーブル・エントリはデコミットされたクラスタをスキップするように修正することができる。デコミットされたスペースがファイルの元の先頭から連続しているときは、メタデータとアプリケーション・プログラムの両方のゼロ・ポイントがシフトされる。ただし、ここで注目すべきことは、このような最適化は、デコミットされたデータがファイルの先頭で連続していないときは、アプリケーション・プログラム30が各不連続の、デコミットされたファイル・セクションごとにそのポインタを調整する必要があるため複雑化することである。

20

詳しく上述した説明から理解されるように、ファイル・システムにおいてファイルの任意の論理部分に割り振られたディスク・スペースを解放する方法とメカニズムが提供されている。この方法とメカニズムはディスク・スペースの解放を高速化すると共に、大量のデータをコピーしたり、大量の一時スペースを割り振ったりしないで済むようにする。この方法とメカニズムはファイル・システムに組み込まれ、ファイル・システムが具備する既存セーフガードと一緒に機能してシステム障害から保護し、マージ・アプリケーションや持続的FIFOキューを取り扱うアプリケーションなどの、アプリケーション・プログラムで容易に利用することができる。この方法とメカニズムは高速、単純で、信頼性があり、拡張可能であり、ほとんどどのファイル・システムでも有効に働く。

30

本発明は種々態様に変更し、代替構成も可能であるが、上記ではいくつかの実施例を図面に示して、詳しく説明してきた。しかし、本発明は上述した特定実施形態に限定されるものではなく、その中には本発明の精神と範囲に属する、すべての変更や改良、代替構成、および等価技術が含まれるものである。

【 図 1 】

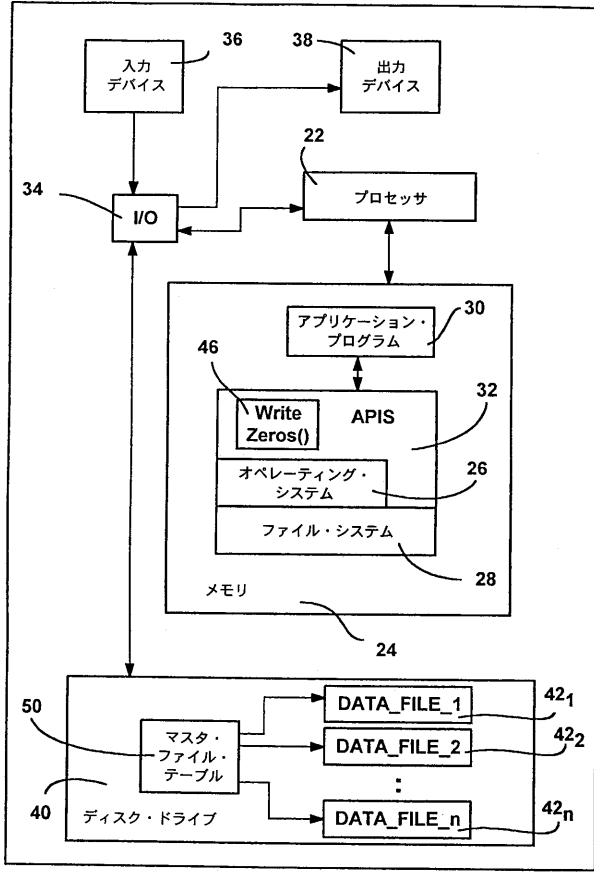


FIG. 1

【 図 2 】

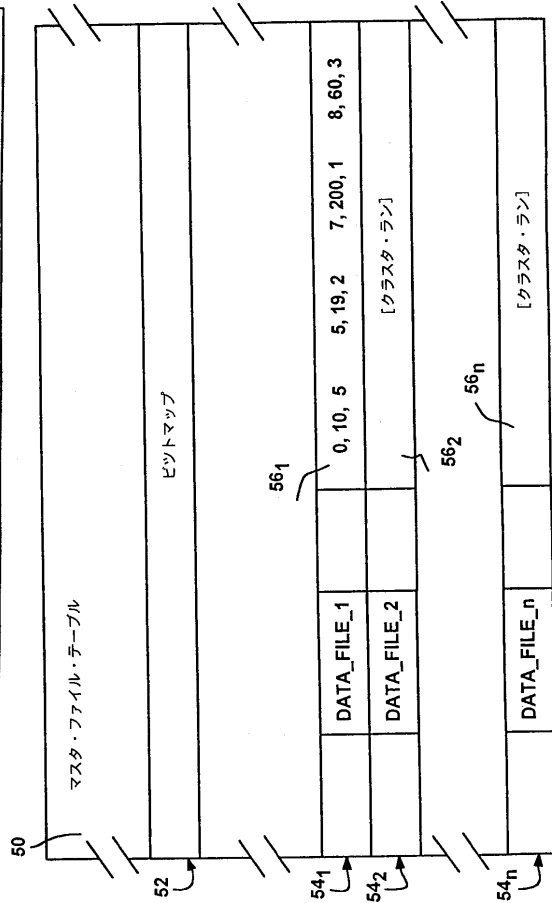


FIG. 2

【 図 3 】

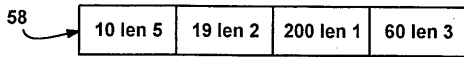


FIG. 3

【 図 4 】

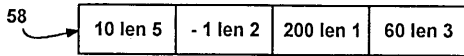


FIG. 4

【 図 5 】

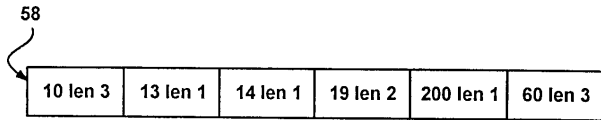


FIG. 5

【 図 6 】

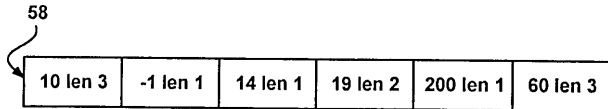


FIG. 6

【 図 8 】

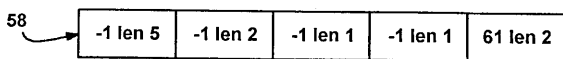


FIG. 8

【 図 9 】

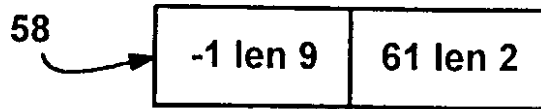
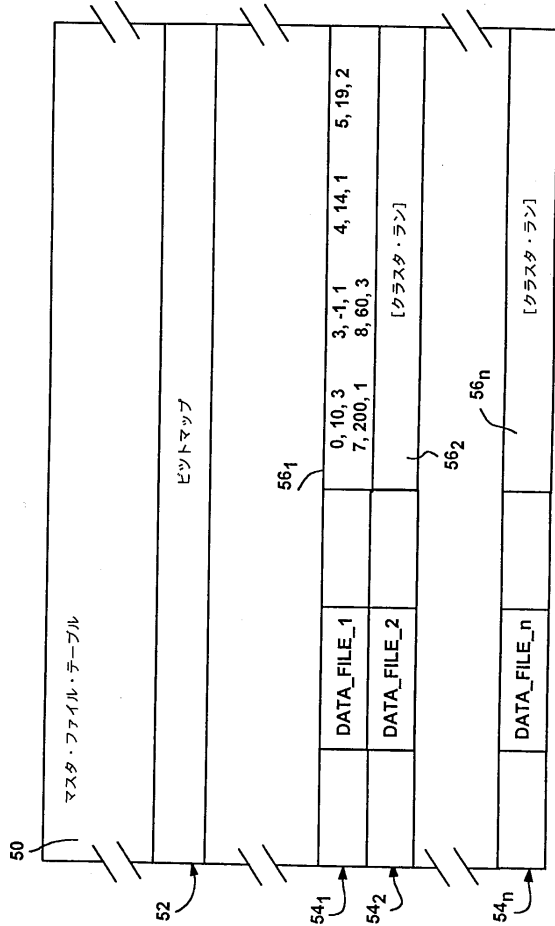


FIG. 9

【 図 7 】



【 図 1 2 】

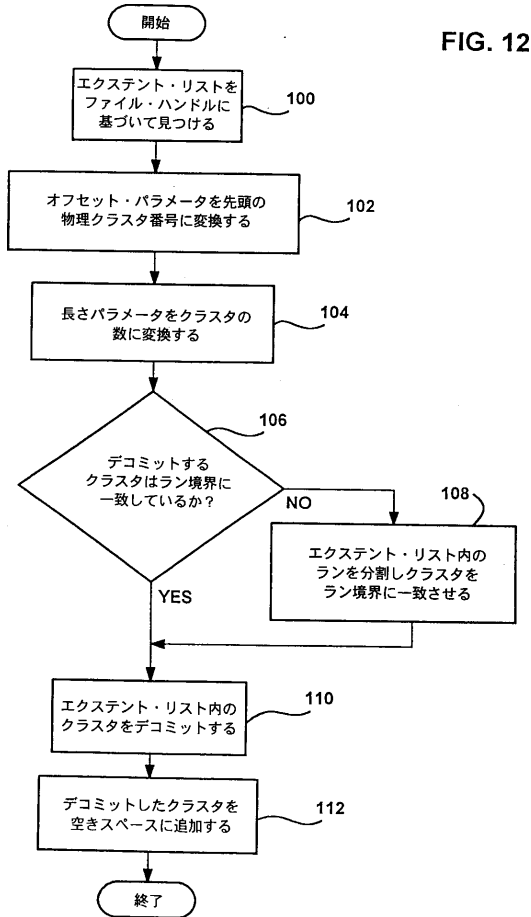


FIG. 12

【 図 1 0 】

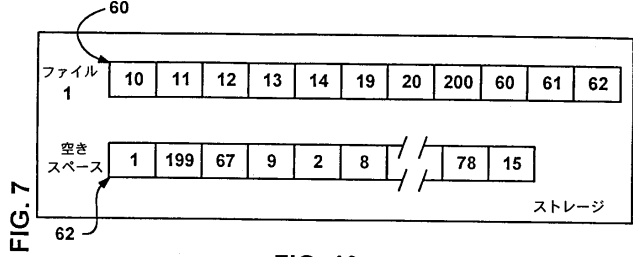


FIG. 7

FIG. 10

【 図 1 1 】

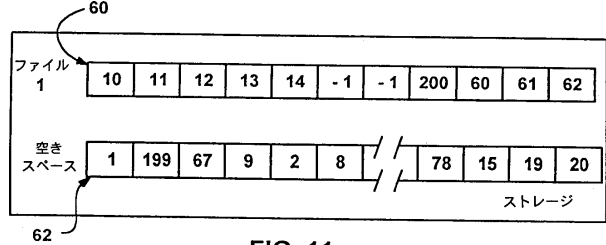


FIG. 11

【 図 1 3 】

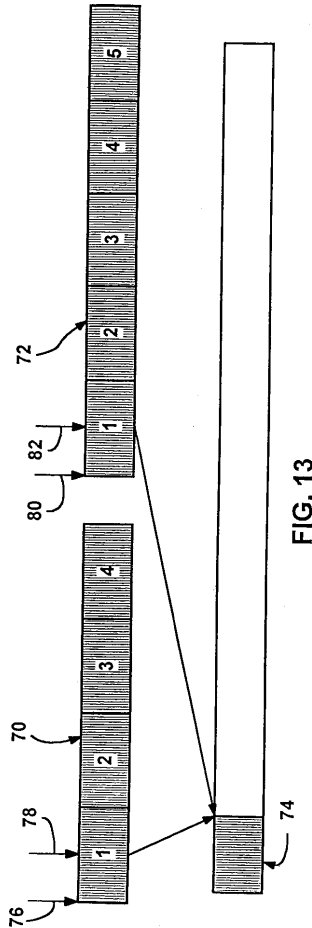
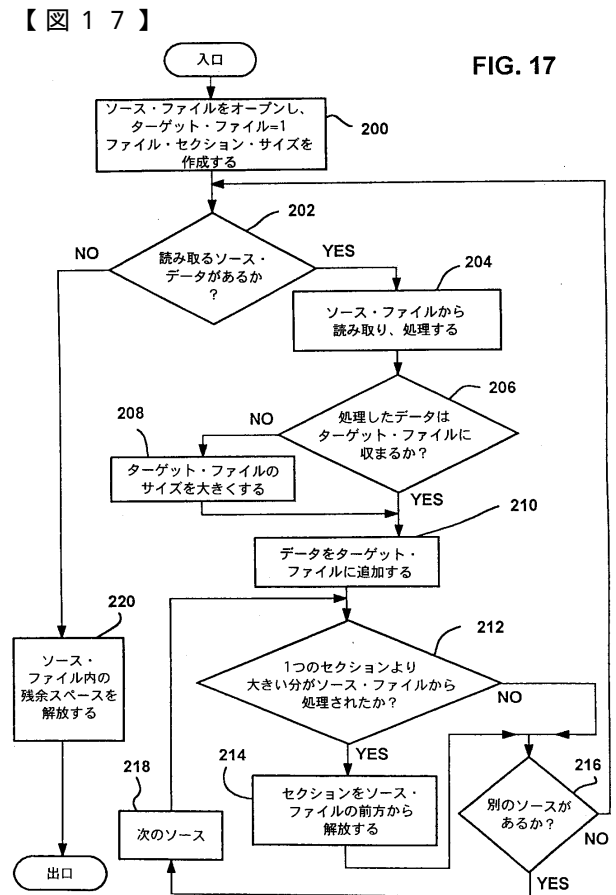
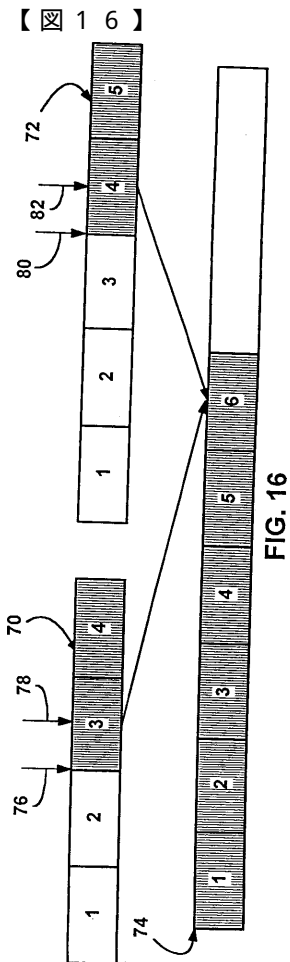
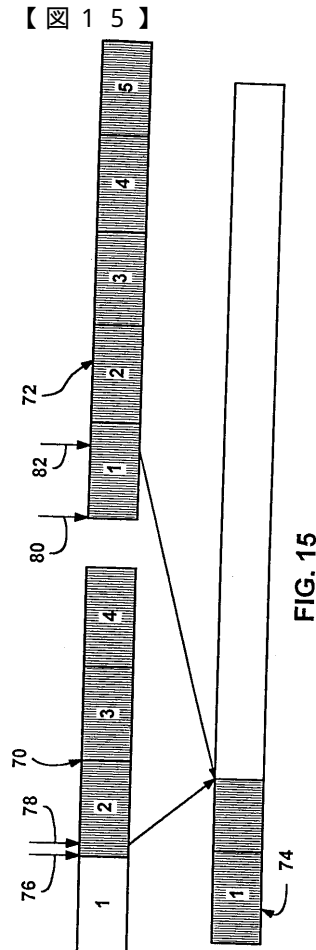
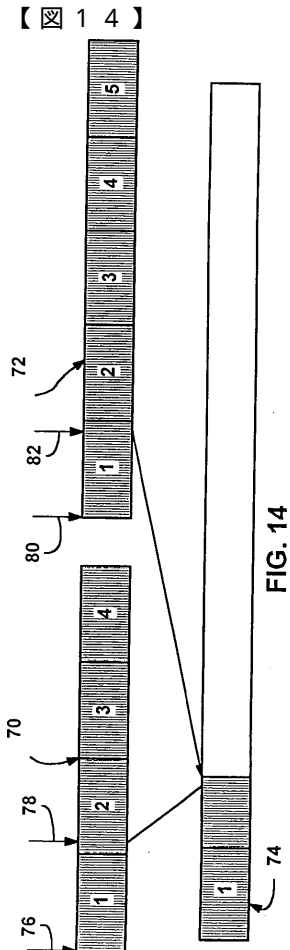


FIG. 13



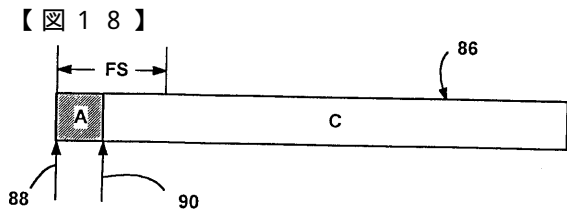


FIG. 18

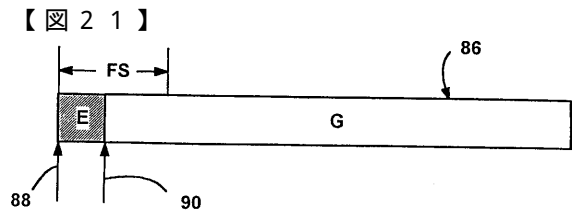


FIG. 21

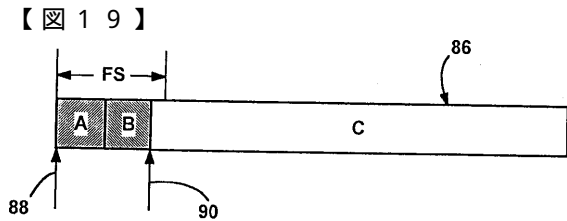


FIG. 19

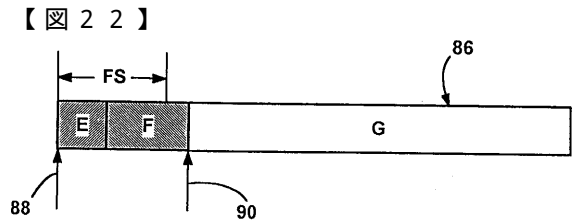


FIG. 22

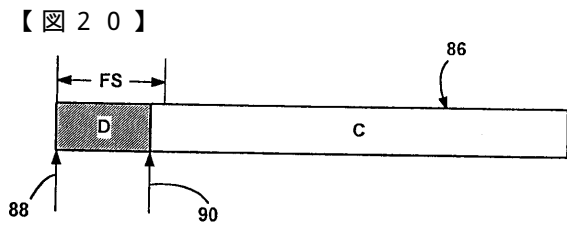


FIG. 20

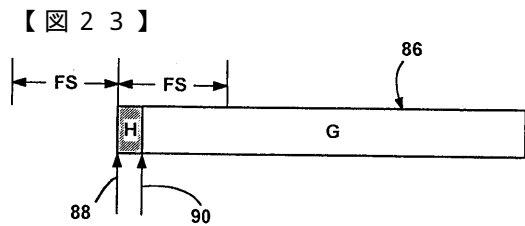
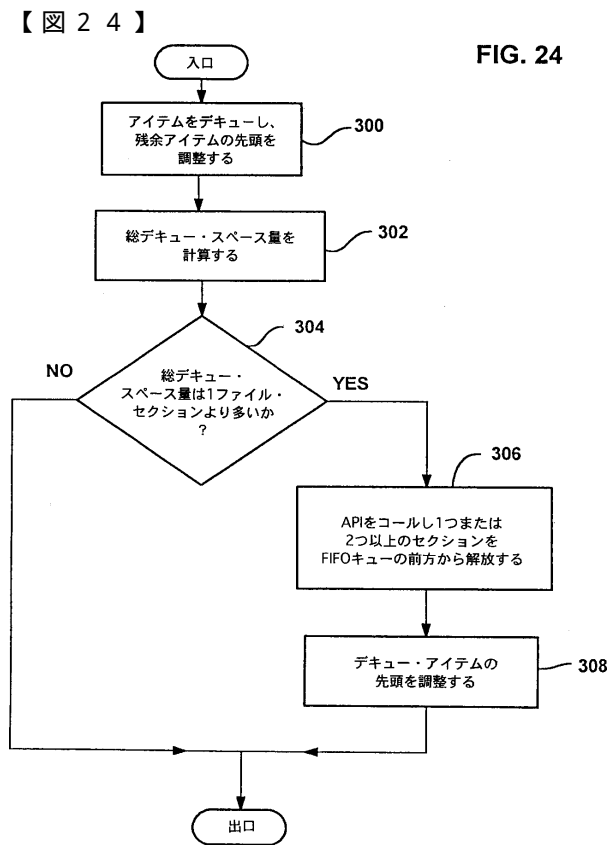


FIG. 23



フロントページの続き

- (72)発明者 ペルトネン, カイル, ジー.
アメリカ合衆国 98027 ワシントン州 アイザッカー サウスイースト 42エヌディー
プレイス 18126
- (72)発明者 バーコウィッツ, ブライアン, ティー.
アメリカ合衆国 98007 ワシントン州 ベルヴィー 142エヌディー プレイス ノース
イースト 3912
- (72)発明者 ズビコウスキー, マーク, ジェイ.
アメリカ合衆国 98072 ワシントン州 ウッディンヴィル ノースイースト 178ティ
エイチ プレイス 15817
- (72)発明者 マイルスキー, パートズ, ビー.
アメリカ合衆国 98125 ワシントン州 シアトル46 ティーエイチ アヴェニュー ノース
イースト 10052

審査官 原 秀人

- (56)参考文献 特開平08-328970(JP, A)
国際公開第96/18960(WO, A1)
Maurice J. Bach, UNIXカーネルの設計, 日本, 共立出版株式会社, 1991年 6月10
日, 第1版, p. 95 - 102
A. V. エイホ 外2名, 情報処理シリーズ11 データ構造とアルゴリズム, 日本, 培風社,
1987年 3月10日, 第1版, p. 48 - 53
MVS/ESA アプリケーション開発解説書: アセンブラー呼出し可能サービス(Open E
dition MVS用) MVS/ESAシステム・プロダクト, 日本アイ・ビー・エム株式会
社, 1996年 5月, 第5版, p. 170-172, 735-738, 1059-1060
Samuel J. Leffer 外3名, The Design and Implementation of the 4.3BSD UNIX Operating Sys
tem, 米国, Addison-Wesley Publishing Company, Inc., 1989年, p. 29 - 34

(58)調査した分野(Int.Cl., DB名)

G06F 12/00